


## TCSS 422: OPERATING SYSTEMS

### CPU Scheduling, Multi-level Feedback Queue (MLFQ)

**Wes J. Lloyd**  
 School of Engineering and Technology,  
 University of Washington - Tacoma



October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
Tacoma

## FEEDBACK FROM 10/15

- **ASIDE: Real-time operating systems:** From: Real time computing Wikipedia: OS for real-time applications that process data as it comes in, typically without buffer delays.
- Assume fixed # of processes; control delay; avoid latency
- Use Case: real-time audio recording/editing
- Types:
  - **Hard:** Missing a deadline results in total system failure
  - **Firm:** Can tolerate infrequent deadline misses, degrades QoS
  - **Soft:** Usefulness of results degrades after deadlines, results in Quality of Service (QoS) degradation
- Linux: "Soft" support via low latency kernel patch
- RTLinux: hard realtime OS microkernel, runs Linux kernel

October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.2

## FEEDBACK - 2

- Real-time job priority in Linux:
- Linux Completely Fair Scheduler:
- Jobs scheduled with a **Real-time policy:**
- SCHED\_FIFO (FF), SCHED\_RR (RR)
  
- Not scheduled as a typical task

October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.3

## FEEDBACK - 3

- What is a (CPU) core?
- How many do most computers have?
- What can you do with them related to threads?

October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.4

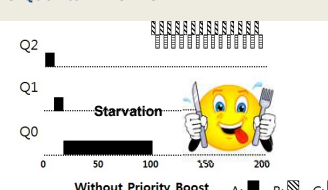
## MLFQ REVIEW

- Without priority boost:
- **Rule 1:** If Priority(A) > Priority(B), A runs (B doesn't).
- **Rule 2:** If Priority(A) = Priority(B), A & B run in RR.
  
- **KEY:** If time quantum of a higher queue is filled, then we don't run any jobs in lower priority queues!!!

October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.5

## STARVATION EXAMPLE

- **Consider 3 queues:**
- Q2 - HIGH PRIORITY - Time Quantum 10ms
- Q1 - MEDIUM PRIORITY - Time Quantum 20 ms
- Q0 - LOW PRIORITY - Time Quantum 40 ms
  
- Job A: 200ms no I/O
- Job B: 5ms then I/O
- Job C: 5ms then I/O
- Q2 fills, starves Q1 & Q0
- A makes no progress



October 15, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.6

### EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms in its highest queue, how often would you have to boost jobs back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?
- Some combination of n short jobs runs for a total of 10 ms per cycle without relinquishing the CPU
  - E.g. 2 jobs = 5 ms ea; 3 jobs = 3.33 ms ea, 10 jobs = 1 ms ea
  - n jobs always uses full time quantum (10 ms)
  - Batch jobs starts, runs for full quantum of 10ms
  - All other jobs run and context switch totaling the quantum per cycle
  - If 10ms is 5% of the CPU, when must the priority boost be ???
  - Priority boost occurs at every 200ms


October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.7

### OBJECTIVES

- C Tutorial (Sunday 10/21)
- Program 1 – MASH Shell (Friday 10/26)
- **CPU Scheduling cont'd:**
- Chapter 8 – Multi-level Feedback Queue
- Chapter 9 – Proportional Share Scheduler
- Linux - Completely Fair Scheduler (CFS)
- **Multi-threaded Programming**
- Chapter 26 – Concurrency Introduction
- Chapter 27 – Linux Thread API
- Chapter 28 – Introduction to Locks

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.8

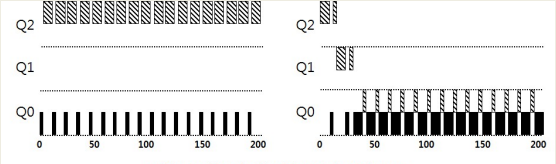
## CHAPTER 8 – MULTI-LEVEL FEEDBACK QUEUE (MLFQ) SCHEDULER



October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.9

### PREVENTING GAMING

- Improved time accounting:
  - Track total job execution time in the queue
  - Each job receives a fixed time allotment
  - When allotment is exhausted, job priority is lowered

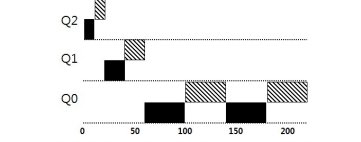


Without(Left) and With(Right) Gaming Tolerance

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.10

### MLFQ: TUNING

- Consider the tradeoffs:
  - How many queues?
  - What is a good time slice?
  - How often should we “Boost” priority of jobs?
  - What about different time slices to different queues?



Example) 10ms for the highest queue, 20ms for the middle, 40ms for the lowest

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.11

### MLFQ RULE SUMMARY

- The refined set of MLFQ rules:
  - **Rule 1:** If Priority(A) > Priority(B), A runs (B doesn't).
  - **Rule 2:** If Priority(A) = Priority(B), A & B run in RR.
  - **Rule 3:** When a job enters the system, it is placed at the highest priority.
  - **Rule 4:** Once a job uses up its time allotment at a given level (regardless of how many times it has given up the CPU), its priority is reduced(i.e., it moves down on queue).
  - **Rule 5:** After some time period S, move all the jobs in the system to the topmost queue.

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.12

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

Job	Arrival Time	Job Length
A	T=0	4
B	T=0	16
C	T=0	8

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will lose points.

## CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.14

### PROPORTIONAL SHARE SCHEDULER

- Also called fair-share scheduler or lottery scheduler
  - Guarantees each job receives some percentage of CPU time based on share of "tickets"
  - Each job receives an allotment of tickets
  - % of tickets corresponds to potential share of a resource
  - Can conceptually schedule any resource this way
    - CPU, disk I/O, memory

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.15

### LOTTERY SCHEDULER

- Simple implementation
  - Just need a random number generator
    - Picks the winning ticket
  - Maintain a data structure of jobs and tickets (list)
  - Traverse list to find the owner of the ticket
  - Consider sorting the list for speed

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.16

### LOTTERY SCHEDULER IMPLEMENTATION

```

1 // counter: used to track if we've found the winner yet
2 int counter = 0;
3
4 // winner: use some call to a random number generator to
5 // get a value, between 0 and the total # of tickets
6 int winner = getRandom(0, totaltickets);
7
8 // current: use this to walk through the list of jobs
9 node* current = head;
10
11 // loop until the sum of ticket values is > the winner
12 while (current) {
13     counter = counter + current->tickets;
14     if (counter > winner)
15         break; // found the winner
16     current = current->next;
17 }
18 // 'current' is the winner: schedule it...
    
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.17

### TICKET MECHANISMS

- Ticket currency / exchange
  - User allocates tickets in any desired way
  - OS converts user currency into global currency
- Example:
  - There are 200 global tickets assigned by the OS

User A → 500 (A's currency) to A1 → 50 (global currency)  
 → 500 (A's currency) to A2 → 50 (global currency)

User B → 10 (B's currency) to B1 → 100 (global currency)

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.18

## TICKET MECHANISMS - 2

- Ticket transfer
  - Temporarily hand off tickets to another process
- Ticket inflation
  - Process can temporarily raise or lower the number of tickets it owns
  - If a process needs more CPU time, it can boost tickets.

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.19

## LOTTERY SCHEDULING

- Scheduler picks a **winning** ticket
  - Load the job with the winning ticket and run it
- Example:
  - Given 100 tickets in the pool
    - Job A has 75 tickets: 0 - 74
    - Job B has 25 tickets: 75 - 99

Scheduler's winning tickets: 63 85 70 39 76 17 29 41 36 39 10 99 68 83 63

Scheduled job: A B A A B A A A A A A B A B A

- But what do we know about probability of a coin flip?

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.20

## COIN FLIPPING

- Equality of distribution (fairness) requires a lot of flips!

Similarly, Lottery scheduling requires lots of "rounds" to achieve fairness.

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.21

## LOTTERY FAIRNESS

- With two jobs
  - Each with the same number of tickets (t=100)

When the job length is not very long, average unfairness can be quite severe.

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.22

## LOTTERY SCHEDULING CHALLENGES

- What is the best approach to assign tickets to jobs?
  - Typical approach is to assume users know best
  - Users are provided with tickets, which they allocate as desired
- How should the OS automatically distribute tickets upon job arrival?
  - What do we know about incoming jobs a priori ?
  - Ticket assignment is really an open problem...

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.23

## STRIDE SCHEDULER

- Addresses statistical probability issues with lottery scheduling
- Instead of guessing a random number to select a job, simply count...

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.24

### STRIDE SCHEDULER - 2

- Jobs have a "stride" value
  - A stride value describes the counter pace when the job should give up the CPU
  - Stride value is **inverse in proportion** to the job's number of tickets (more tickets = smaller stride)
- Total system tickets = 10,000
  - Job A has 100 tickets →  $A_{stride} = 10000/100 = 100$  stride
  - Job B has 50 tickets →  $B_{stride} = 10000/50 = 200$  stride
  - Job C has 250 tickets →  $C_{stride} = 10000/250 = 40$  stride
- Stride scheduler tracks "pass" values for each job (A, B, C)

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.25

### STRIDE SCHEDULER - 3

- Basic algorithm:
  - Stride scheduler picks job with the lowest pass value
  - Scheduler increments job's pass value by its stride and starts running
  - Stride scheduler increments a counter
  - When counter exceeds pass value of current job, pick a new job (go to 1)
- KEY:** When the counter reaches a job's "PASS" value, the scheduler passes on to the next job...

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.26

### STRIDE SCHEDULER - EXAMPLE

- Stride values
  - Tickets = priority to select job
  - Stride is inverse to tickets
  - Lower stride = more chances to run (higher priority)

**Priority**  
 C stride = 40  
 A stride = 100  
 B stride = 200

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.27

### STRIDE SCHEDULER EXAMPLE - 2

- Three-way tie: randomly pick job A (all pass values=0)
- Set A's pass value to A's stride = 100
- Increment counter until > 100
- Pick a new job: two-way tie

Pass(A) (stride=100)	Pass(B) (stride=200)	Pass(C) (stride=40)	Who Runs?
0	0	0	A
100	0	0	B
100	200	0	C
100	200	40	C
100	200	80	C
100	200	120	A
200	200	120	C
200	200	160	C
200	200	200	...

**Tickets**  
 C = 250  
 A = 100  
 B = 50

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.28

### STRIDE SCHEDULER EXAMPLE - 3

- We set A's counter (pass value) to A's stride = 100
- Next scheduling decision between B (pass=0) and C (pass=0)
  - Randomly choose B
- C has the lowest counter for next 3 rounds

Pass(A) (stride=100)	Pass(B) (stride=200)	Pass(C) (stride=40)	Who Runs?
0	0	0	A
100	0	0	B
100	200	0	C
100	200	40	C
100	200	80	C
100	200	120	A
200	200	120	C
200	200	160	C
200	200	200	...

**Tickets**  
 C = 250  
 A = 100  
 B = 50

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.29

### STRIDE SCHEDULER EXAMPLE - 4

- Job counters support determining which job to run next
- Over time jobs are scheduled to run based on their priority represented as their **share of tickets...**
- Tickets are analogous to job priority**

Pass(A) (stride=100)	Pass(B) (stride=200)	Pass(C) (stride=40)	Who Runs?
0	0	0	A
100	0	0	B
100	200	0	C
100	200	40	C
100	200	80	C
100	200	120	A
200	200	120	C
200	200	160	C
200	200	200	...

**Tickets**  
 C = 250  
 A = 100  
 B = 50

October 15, 2018    TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.30

### LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Linux ≥ 2.6.23: Completely Fair Scheduler (CFS)
- Linux < 2.6.23: O(1) scheduler
- Every thread/process has a scheduling policy:
- **Normal policies:** SCHED\_OTHER (TS), SCHED\_IDLE, SCHED\_BATCH
  - TS = Time Sharing
- **Real-time policies:** SCHED\_FIFO (FF), SCHED\_RR (RR)
- Show scheduling policy and priority:
  - `ps -elfc`
  - `ps ax -o pid,ni,cls,pri,cmd`

October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.31
------------------	---	-------

### COMPLETELY FAIR SCHEDULER - 2

- Loosely based on the stride scheduler
- CFS models system as a Perfect Multi-Tasking System
  - In perfect system every process of the same priority receives exactly  $1/n^{\text{th}}$  of the CPU time
- Scheduling classes (runqueues)
  - Groups process of same priority across set of runqueues
  - Process priority groups use different sets of runqueues for priorities
    - Default (SCHED\_OTHER) gets a set (PRI 1-99)
    - Real-time (FF,RR) separate sets (PRI 1-139)
  - Scheduler picks task with lowest accumulative runtime to run
  - Time quantum based on proportion of CPU time (%), not fixed time allotments
  - Quantum varies based on how many jobs in shared runqueue

October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.32
------------------	---	-------

### COMPLETELY FAIR SCHEDULER - 3

- CFS uses weighted fair queuing
  - 1<sup>st</sup> implementation of fair queuing process scheduler in a major OS
- Runqueues are stored using a linux red-black tree
  - Self balancing binary search tree- nodes indexed by `vruntime`
  - Leftmost node has lowest `vruntime` (total execution time)
  - Walking tree to find left most node is only  $O(\log N)$  for N nodes
  - Completed processes removed
  - Processes using up quantum, or interrupted reinserted  
*They are in READY state...*
  - This way, processes that sleep a lot (i.e. event handlers) have low `vruntime`, get a "priority boost" when they need to run
- Key takeaway  
identifying the next job to schedule is **really** fast!

October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.33
------------------	---	-------

### COMPLETELY FAIR SCHEDULER - 4

- Time slice: Linux **"Nice value"**
  - Nice value predates the CFS scheduler
  - Top shows nice values
  - Process command (nice & priority):  
`ps ax -o pid,ni,cmd,%cpu, pri`
- Nice Values: from -20 to 19
  - Lower is **higher** priority, default is 0
  - Scheduling quantum is calculated using nice value
  - Default: `cat /proc/sys/kernel/sched_rr_timeslice_ms`
  - Target latency:
    - Interval during which task should run at least once
    - Automatically increases as number of jobs increases

October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.34
------------------	---	-------

### COMPLETELY FAIR SCHEDULER - 5

- Challenge:
  - How do we map a nice value to an actual CPU timeslice (ms)?
  - What is the best mapping?
    - O(1) scheduler (< 2.6.23)
      - tried to map nice value to timeslice (fixed allotment)
  - Linux completely fair scheduler
    - Nice value suggests priority used to assign runqueue for job
    - Time proportion varies based on # of jobs in runqueue
      - with fewer jobs in runqueue, time proportion is larger

October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.35
------------------	---	-------

### COMPLETELY FAIR SCHEDULER - 6

- Nice values become relative for determining time slices
  - Proportion of CPU time to allocate is relative to other queued tasks
- Scheduler tracks virtual run time in `vruntime` variable
- The task on a given runqueue (nice value) with the lowest `vruntime` is scheduled next
- `struct sched_entity` contains `vruntime` parameter
  - Describes process execution time in nanoseconds
- Perfect scheduler →  
achieve equal `vruntime` for all processes of same priority


October 15, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.36
------------------	---	-------

## COMPLETELY FAIR SCHEDULER - 7

- More information:
- Man page: "man sched" : Describes Linux scheduling API
- <http://manpages.ubuntu.com/manpages/bionic/man7/sched.7.html>
- <https://www.kernel.org/doc/Documentation/scheduler/sched-design-CFS.txt>
- [https://en.wikipedia.org/wiki/Completely\\_Fair\\_Scheduler](https://en.wikipedia.org/wiki/Completely_Fair_Scheduler)
- See paper: The Linux Scheduler – a Decade of Wasted Cores
- <http://www.ece.ubc.ca/~sasha/papers/eurosys16-final29.pdf>

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.37

## CHAPTER 26 - CONCURRENCY: AN INTRODUCTION



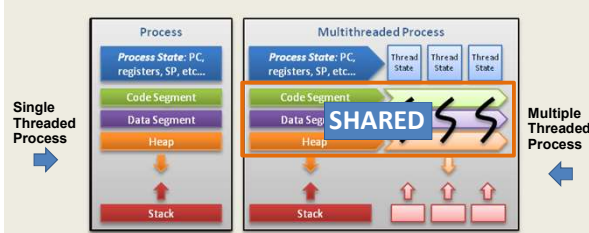
October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.38

## OBJECTIVES

- Introduction to threads
- Race condition
- Critical section
- Thread API

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.39

## THREADS



©Alfred Park, <http://randu.org/tutorials/threads>

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.40

## THREADS - 2

- Enables a single process (program) to have multiple "workers"
- Supports independent path(s) of execution within a program *with shared memory* ...
- Each thread has its own Thread Control Block (TCB)
  - PC, registers, SP, and stack
- Code segment, memory, and heap are shared

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.41

## PROCESS AND THREAD METADATA

- Thread Control Block vs. Process Control Block

**Thread Identification**

Thread state

CPU information:

- Program counter
- Register contents

Thread priority

Pointer to process that created this thread

Pointers to all other threads created by this thread

**Process Identification**

Process status

Process state:

- Process status word
- Register contents
- Main memory
- Resources
- Process priority

Accounting

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.42

## SHARED ADDRESS SPACE

- Every thread has it's own stack / PC

**A Single-Threaded Address Space**

The code segment: where instructions live

The heap segment: contains malloc'd data dynamic data structures (it grows downward)

(it grows upward)

The stack segment: contains local variables arguments to routines, return values, etc.

**Two threaded Address Space**

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.43

## THREAD CREATION EXAMPLE

```

#include <stdio.h>
#include <assert.h>
#include <pthread.h>

void *mythread(void *arg) {
    printf("%s\n", (char *) arg);
    return NULL;
}

int
main(int argc, char *argv[]) {
    pthread_t p1, p2;
    int rc;
    printf("main: begin\n");
    rc = pthread_create(&p1, NULL, mythread, "A"); assert(rc == 0);
    rc = pthread_create(&p2, NULL, mythread, "B"); assert(rc == 0);
    // join waits for the threads to finish
    rc = pthread_join(p1, NULL); assert(rc == 0);
    rc = pthread_join(p2, NULL); assert(rc == 0);
    printf("main: end\n");
    return 0;
}
    
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.44

## POSSIBLE ORDERINGS OF EVENTS

int main()	Thread 1	Thread 2
Starts running		
Prints 'main: begin'		
Creates Thread 1		
Creates Thread 2		
Waits for T1		
	Runs	
	Prints 'A'	
	Returns	
Waits for T2		
		Runs
		Prints 'B'
		Returns
Prints 'main: end'		

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.45

## POSSIBLE ORDERINGS OF EVENTS - 2

int main()	Thread 1	Thread 2
Starts running		
Prints 'main: begin'		
Creates Thread 1		
	Runs	
	Prints 'A'	
	Returns	
Creates Thread 2		
		Runs
		Prints 'B'
		Returns
Waits for T1	Returns immediately	
Waits for T2		Returns immediately
Prints 'main: end'		

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.46

## POSSIBLE ORDERINGS OF EVENTS - 3

int main()	Thread 1	Thread 2
Starts running		
Prints 'main: begin'		
Creates Thread 1		
Creates Thread 2		
Waits for T1		
	Runs	
	Prints 'A'	
	Returns	
Waits for T2		
		Immediately returns
Prints 'main: end'		

What if execution order of events in the program matters?

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.47


## COUNTER EXAMPLE

- Counter example
- A + B : ordering
- Counter: incrementing global variable by two threads

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.48



## PROCESSES VS. THREADS



- What's the difference between forks and threads?
  - Forks:** duplicate a process
  - Think of **CLONING** - There will be two identical processes at the end
  - Threads:** no duplicate of code/heap, lightweight execution threads

**Process**

Process State: PC, registers, SP, etc...

Code Segment

Data Segment

Heap

Stack

**Process**

Process State: PC, registers, SP, etc...

Code Segment

Data Segment

Heap

Stack

code data files

registers stack

thread

single-threaded process

code data files

registers registers registers

stack stack stack

thread

multithreaded process

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.49

## RACE CONDITION

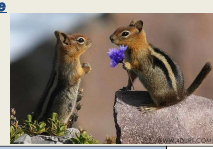
- What is happening with our counter?
  - When counter=50, consider code: counter = counter + 1
  - If synchronized, counter will = 52

	OS	Thread1	Thread2	(after instruction)	PC	%eax	counter
		before critical section			100	0	50
		mov 0x8049a1c, %eax			105	50	50
		add \$0x1, %eax			108	51	50
Interrupt		save T1's state			100	0	50
		restore T2's state			105	50	50
			mov 0x8049a1c, %eax		108	51	50
			add \$0x1, %eax		108	51	50
			mov %eax, 0x8049a1c		113	51	51
Interrupt		save T2's state			108	51	50
		restore T1's state			108	51	50
		mov %eax, 0x8049a1c			113	51	51

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.50

## CRITICAL SECTION

- Code that accesses a shared variable must not be **concurrently** executed by more than one thread
- Multiple active threads inside a **critical section** produce a **race condition**.
- Atomic execution** (all code executed as a *unit*) must be ensured in **critical** sections
  - These sections must be **mutually exclusive**



October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.51

## LOCKS

- To demonstrate how critical section(s) can be executed "atomically-as a *unit*" Chapter 27 & beyond introduce locks

```

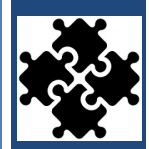
1 lock_t mutex;
2 . . .
3 lock(&mutex);
4 balance = balance + 1;
5 unlock(&mutex);
    
```

Critical section

- Counter example revisited

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.52

# CHAPTER 27 - LINUX THREAD API



October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.53

## THREAD CREATION

- pthread\_create

```

#include <pthread.h>

int
pthread_create( pthread_t* thread,
               const pthread_attr_t* attr,
               void* (*start_routine)(void*),
               void* arg);
    
```

- thread: thread struct
- attr: stack size, scheduling priority... (optional)
- start\_routine: function pointer to thread routine
- arg: argument to pass to thread routine (optional)

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma
L6.54

### PTHREAD\_CREATE – PASS ANY DATA

```

#include <pthread.h>

typedef struct _myarg_t {
    int a;
    int b;
} myarg_t;

void *mythread(void *arg) {
    myarg_t *m = (myarg_t *) arg;
    printf("%d %d\n", m->a, m->b);
    return NULL;
}

int main(int argc, char *argv[]) {
    pthread_t p;
    int rc;

    myarg_t args;
    args.a = 10;
    args.b = 20;
    rc = pthread_create(&p, NULL, mythread, &args);
}
    
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.55

### PASSING A SINGLE VALUE

Using this approach on your Ubuntu VM, How large (in bytes) can the primitive data type be?

```

9 int rc, m;
10 pthread_create(&p, NULL, mythread, (void *)100);
11 pthread_join(p, (void **) &m);
12 printf("returned %d\n", m);
13 return 0;
14 )
    
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.56

### WAITING FOR THREADS TO FINISH

```
int pthread_join(pthread_t thread, void **value_ptr);
```

- thread: which thread?
- value\_ptr: pointer to return value  
type is dynamic / agnostic
- Returned values *must* be on the heap
- Thread stacks destroyed upon thread termination (join)
- Pointers to thread stack memory addresses are invalid
  - May appear as gibberish or lead to crash (seg fault)
- Not all threads join – *What would be Examples ??*

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.57

### What will this code do?

```

struct myarg {
    int a;
    int b;
};

void *worker(void *arg)
{
    struct myarg *input = (struct myarg *) arg;
    printf("a=%d b=%d\n", input->a, input->b);
    struct myarg output;
    output.a = 1;
    output.b = 2;
    return (void *) &output;
}

int main (int argc, char * argv[])
{
    pthread_t p1;
    struct myarg args;
    struct myarg *ret_args;
    args.a = 10;
    args.b = 20;
    pthread_t p;
    pthread_create(&p1, NULL, worker, &args);
    pthread_join(p1, (void **) &ret_args);
    printf("returned %d %d\n", ret_args->a, ret_args->b);
    return 0;
}
    
```

➔ Data on thread stack

```

$ ./pthread_struct
a=10 b=20
Segmentation fault (core dumped)
    
```

### How can this code be fixed?

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.58

### How about this code?

```

struct myarg {
    int a;
    int b;
};

void *worker(void *arg)
{
    struct myarg *input = (struct myarg *) arg;
    printf("a=%d b=%d\n", input->a, input->b);
    input->a = 1;
    input->b = 2;
    return (void *) &input;
}

int main (int argc, char * argv[])
{
    pthread_t p1;
    struct myarg args;
    struct myarg *ret_args;
    args.a = 10;
    args.b = 20;
    pthread_create(&p1, NULL, worker, &args);
    pthread_join(p1, (void **) &ret_args);
    printf("returned %d %d\n", ret_args->a, ret_args->b);
    return 0;
}
    
```

```

$ ./pthread_struct
a=10 b=20
returned 1 2
    
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.59

### ADDING CASTS

- Casting
- Suppresses compiler warnings when passing "typed" data where (void) or (void \*) is called for
- Example: uncasted capture in pthread\_join  
pthread\_int.c: In function 'main':  
pthread\_int.c:34:20: warning: passing argument 2 of 'pthread\_join' from incompatible pointer type [-Wincompatible-pointer-types]  
pthread\_join(p1, &ival);
- Example: uncasted return  
In file included from pthread\_int.c:3:0:  
/usr/include/pthread.h:250:12: note: expected 'void \*\*' but argument is of type 'int \*\*'  
extern int pthread\_join (pthread\_t \_\_th, void \*\*\_\_thread\_return);

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.60

## ADDING CASTS - 2

- pthread\_join
 

```
int * p1val;
int * p2val;
pthread_join(p1, (void *)&p1val);
pthread_join(p2, (void *)&p2val);
```
- return from thread function
 

```
int * counterval = malloc(sizeof(int));
*counterval = counter;
return (void *) counterval;
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.61

## LOCKS

- pthread\_mutex\_t data type
- /usr/include/bits/pthread\_types.h
 

```
// Global Address Space
static volatile int counter = 0;
pthread_mutex_t lock;

void *worker(void *arg)
{
    int i;
    for (i=0;i<1000000;i++) {
        int rc = pthread_mutex_lock(&lock);
        assert(rc==0);
        counter = counter + 1;
        pthread_mutex_unlock(&lock);
    }
    return NULL;
}
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.62

## LOCKS - 2

- Ensure critical sections are executed atomically-as a unit
  - Provides implementation of "Mutual Exclusion"
- API
 

```
int pthread_mutex_lock(pthread_mutex_t *mutex);
int pthread_mutex_unlock(pthread_mutex_t *mutex);
```
- Example w/o initialization & error checking
 

```
pthread_mutex_t lock;
pthread_mutex_lock(&lock);
x = x + 1; // or whatever your critical section is
pthread_mutex_unlock(&lock);
```

  - Blocks forever until lock can be obtained
  - Enters critical section once lock is obtained
  - Releases lock

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.63

## LOCK INITIALIZATION

- Assigning the constant
 

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
```
- API call:
 

```
int rc = pthread_mutex_init(&lock, NULL);
assert(rc == 0); // always check success!
```
- Initializes mutex with attributes specified by 2<sup>nd</sup> argument
- If NULL, then default attributes are used
- Upon initialization, the mutex is initialized and unlocked

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.64

## LOCKS - 3


- Error checking wrapper
 

```
// Use this to keep your code clean but check for failures
// Only use if exiting program is OK upon failure
void Pthread_mutex_lock(pthread_mutex_t *mutex) {
    int rc = pthread_mutex_lock(mutex);
    assert(rc == 0);
}
```
- What if lock can't be obtained?
 

```
int pthread_mutex_trylock(pthread_mutex_t *mutex);
int pthread_mutex_timelock(pthread_mutex_t *mutex,
                           struct timespec *abs_timeout);
```
- trylock – returns immediately (fails) if lock is unavailable
- timelock – tries to obtain a lock for a specified duration

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.65

## CONDITIONS AND SIGNALS

- Condition variables support "signaling" between threads
 
- pthread\_cond\_wait(pthread\_cond\_t \*cond, pthread\_mutex\_t \*mutex);
 

```
int pthread_cond_signal(pthread_cond_t *cond);
```
- pthread\_cond\_t datatype
- pthread\_cond\_wait()
  - Puts thread to "sleep" (waits) (THREAD IS BLOCKED)
  - Threads added to FIFO queue, lock is released
  - Waits (*listens*) for a "signal" (NON-BUSY WAITING, no polling)
  - When signal occurs, interrupt fires, wakes up first thread, (THREAD IS RUNNING), lock is provided to thread

October 15, 2018 TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma L6.66

### CONDITIONS AND SIGNALS - 2

```
int pthread_cond_signal(pthread_cond_t * cond);
int pthread_cond_broadcast(pthread_cond_t * cond);
```

- pthread\_cond\_signal()
  - Called to send a "signal" to wake-up first thread in FIFO "wait" queue
  - The goal is to unblock a thread to respond to the signal
- pthread\_cond\_broadcast()
  - Unlocks all threads in FIFO "wait" queue, currently blocked on the specified condition variable
  - Broadcast is used when all threads should wake-up for the signal
- Which thread is unblocked first?
  - Determined by OS scheduler (based on priority)
  - Thread(s) awoken based on placement order in FIFO wait queue
  - When awoken threads acquire lock as in pthread\_mutex\_lock()

October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.67

### CONDITIONS AND SIGNALS - 3

- Wait example:
 

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

pthread_mutex_lock(&lock);
while (initialized == 0)
    pthread_cond_wait(&cond, &lock);
// Perform work that requires lock
a = a + b;
pthread_mutex_unlock(&lock);
```
- wait puts thread to sleep, releases lock
- when awoken, lock reacquired (but then released by this code)
- When initialized, another thread signals
 

```
pthread_mutex_lock(&lock);
initialized = 1;
pthread_cond_signal(&init);
pthread_mutex_unlock(&lock);
```

State variable set, Enables other thread(s) to proceed above.

October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.68

### CONDITION AND SIGNALS - 4

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

pthread_mutex_lock(&lock);
while (initialized == 0)
    pthread_cond_wait(&cond, &lock);
// Perform work that requires lock
a = a + b;
pthread_mutex_unlock(&lock);
```

- Why do we wait inside a while loop?
  - The while ensures upon awakening the condition is rechecked
  - A signal is raised, but the pre-conditions required to proceed may have not been met. **\*\*MUST CHECK STATE VARIABLE\*\***
  - Without checking the state variable the thread may proceed to execute when it should not. (e.g. too early)

October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.69

### PTHREADS LIBRARY

- Compilation
  - gcc -pthread pthread.c -o pthread
  - Requires explicitly linking the library with compiler flag
  - Use makefile to provide compiler arguments
- List of pthread manpages
  - man -k pthread

October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.70

### SAMPLE MAKEFILE

```
CC=gcc
CFLAGS=-pthread -I. -Wall

binaries=pthread pthread_int pthread_lock_cond pthread_struct
all: $(binaries)


pthread_mult: pthread.c pthread_int.c
$(CC) $(CFLAGS) $^ -o $@

clean:
$(RM) -f $(binaries) *.o
```

- Example builds multiple single file programs
  - All target
- pthread\_mult
  - Example if multiple source files should produce a single executable
- clean target


October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.71

## CHAPTER 28 – LOCKS



October 15, 2018    TCSS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma    L6.72

## LOCKS



- Ensure critical section(s) are executed atomically-as a *unit*
  - Only one thread is allowed to execute a critical section at any given time
  - Ensures the code snippets are "mutually exclusive"
- Protect a global counter:
 

```
balance = balance + 1;
```
- A "critical section":
 

```
1 lock_t mutex; // some globally-allocated lock 'mutex'
2 ...
3 lock (&mutex);
4 balance = balance + 1;
5 unlock (&mutex);
```

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.73

## LOCKS - 2

- Lock variables are called "MUTEX"
  - Short for mutual exclusion (that's what they guarantee)
- Lock variables store the state of the lock
- States
  - Locked** (acquired or held)
  - Unlocked** (available or free)
- Only 1 thread can hold a lock

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.74

## LOCKS - 3

- `pthread_mutex_lock(&lock)`
  - Try to acquire lock
  - If lock is free, calling thread will acquire the lock
  - Thread with lock enters critical section
    - Thread "owns" the lock
- No other thread can acquire the lock before the owner releases it.

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.75

## LOCKS - 4

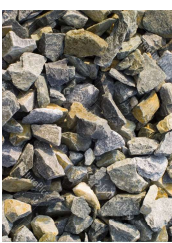
- Program can have many mutex (lock) variables to "serialize" many critical sections
- Locks are also used to protect data structures
  - Prevent multiple threads from changing the same data simultaneously
  - Programmer can make sections of code "granular"
    - Fine grained – means just one grain of sand at a time through an hour glass
  - Similar to relational database transactions
    - DB transactions prevent multiple users from modifying a table, row, field

October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.76

## FINE GRAINED?

- Is this code a good example of "fine grained parallelism"?

```
pthread_mutex_lock(&lock);
a = b++;
b = a * c;
*d = a + b + c;
FILE * fp = fopen ("file.txt", "r");
fscanf(fp, "%s %s %s %d", str1, str2, str3, &e);
ListNode *node = mylist->head;
int i=0
while (node) {
    node->title = str1;
    node->subheading = str2;
    node->desc = str3;
    node->end = *e;
    node = node->next;
    i++
}
e = e - i;
pthread_mutex_unlock(&lock);
```



October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.77

## FINE GRAINED PARALLELISM


```
pthread_mutex_lock(&lock_a);
pthread_mutex_lock(&lock_b);
a = b++;
pthread_mutex_unlock(&lock_b);
pthread_mutex_unlock(&lock_a);

pthread_mutex_lock(&lock_b);
b = a * c;
pthread_mutex_unlock(&lock_b);

pthread_mutex_lock(&lock_d);
*d = a + b + c;
pthread_mutex_unlock(&lock_d);

FILE * fp = fopen ("file.txt", "r");
pthread_mutex_lock(&lock_e);
fscanf(fp, "%s %s %s %d", str1, str2, str3, &e);
pthread_mutex_unlock(&lock_e);


ListNode *node = mylist->head;
int i=0 . . .
```



October 15, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L6.78

## EVALUATING LOCK IMPLEMENTATIONS

- **Correctness**
  - Does the lock work?
  - Are critical sections mutually exclusive? (atomic-as a unit?)
- **Fairness**
  - Are threads competing for a lock have a fair chance of acquiring it?
- **Overhead**



October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.79

## BUILDING LOCKS

- **Locks require hardware support**
  - To minimize overhead, ensure fairness and correctness
  - Special “atomic-as a unit” instructions to support lock implementation
  - Atomic-as a unit exchange instruction
    - XCHG
  - Compare and exchange instruction
    - CMPXCHG
    - CMPXCHG8B
    - CMPXCHG16B

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.80

## HISTORICAL IMPLEMENTATION

- **To implement mutual exclusion**
  - Disable interrupts upon entering critical sections

```


1 void lock() {
2     DisableInterrupts();
3 }
4 void unlock() {
5     EnableInterrupts();
6 }
```

- Any thread could disable system-wide interrupt
  - What if lock is never released?
- On a multiprocessor processor each CPU has its own interrupts
  - Do we disable interrupts for all cores simultaneously?
- While interrupts are disabled, they could be lost
  - If not queued...

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.81

## SPIN LOCK IMPLEMENTATION

- Operate without atomic-as a unit assembly instructions
- “Do-it-yourself” Locks
- Is this lock implementation: Correct? Fair? Performant?



```

1 typedef struct __lock_t { int flag; } lock_t;
2
3 void init(lock_t *mutex) {
4     // 0 → lock is available, 1 → held
5     mutex->flag = 0;
6 }
7
8 void lock(lock_t *mutex) {
9     while (mutex->flag == 1) // TEST the flag
10        ; // spin-wait (do nothing)
11     mutex->flag = 1; // now SET it!
12 }
13
14 void unlock(lock_t *mutex) {
15     mutex->flag = 0;
16 }
```

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.82

## DIY: CORRECT?

- **Correctness requires lock... (e.g. DIY lock is incorrect)**

Thread1	Thread2
call lock() while (flag == 1) interrupt: switch to Thread 2	call lock() while (flag == 1) flag = 1; interrupt: switch to Thread 1
flag = 1; // set flag to 1 (too!)	

- Here both threads have “acquired” the lock simultaneously

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.83

## DIY: PERFORMANT?

```

void lock(lock_t *mutex)
{
    while (mutex->flag == 1); // while lock is unavailable, wait...
    mutex->flag = 1;
}
```

- What is wrong with while(<cond>); ?
- Spin-waiting wastes time actively waiting for another thread
- while (1); will “peg” a CPU core at 100%
  - Continuously loops, and evaluates mutex->flag value...
  - Generates heat...

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.84

## TEST-AND-SET INSTRUCTION

- C implementation: not atomic
  - Adds a simple check to basic spin lock
  - One a single core CPU system with preemptive scheduler:
  - Try this...

```

1 int TestAndSet(int *ptr, int new) {
2     int old = *ptr; // fetch old value at ptr
3     *ptr = new;    // store 'new' into ptr
4     return old;   // return the old value
5 }
```

- lock() method checks that TestAndSet doesn't return 1
- Comparison is in the caller
- Single core systems are becoming scarce
- Try on a one-core VM

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.85

## DIY: TEST-AND-SET - 2

- Requires a preemptive scheduler on single CPU core system
- Lock is never released without a context switch
- 1-core VM: occasionally will deadlock, doesn't miscount

```

1 typedef struct __lock_t {
2     int flag;
3 } lock_t;
4
5 void init(lock_t *lock) {
6     // 0 indicates that lock is available,
7     // 1 that it is held
8     lock->flag = 0;
9 }
10
11 void lock(lock_t *lock) {
12     while (TestAndSet(&lock->flag, 1) == 1)
13         ; // spin-wait
14 }
15
16 void unlock(lock_t *lock) {
17     lock->flag = 0;
18 }
```

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.86

## SPIN LOCK EVALUATION

- **Correctness:**
  - Spin locks guarantee: critical sections won't be executed simultaneously by (2) threads
- **Fairness:**
  - No fairness guarantee. Once a thread has a lock, nothing forces it to relinquish it...
- **Performance:**
  - Spin locks perform "busy waiting"
  - Spin locks are best for short periods of waiting
  - Performance is slow when multiple threads share a CPU
    - Especially for long periods

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.87

## COMPARE AND SWAP

- Checks that the lock variable has the expected value FIRST, before changing its value
  - If so, make assignment
  - Return value at location
- Adds a comparison to TestAndSet
- Useful for wait-free synchronization
  - Supports implementation of shared data structures which can be updated atomically (as a unit) using the HW support CompareAndSwap instruction
  - Shared data structure updates become "wait-free"
  - Upcoming in Chapter 32

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.88

## COMPARE AND SWAP

- Compare and Swap

```

1 int CompareAndSwap(int *ptr, int expected, int new) {
2     int actual = *ptr;
3     if (actual == expected)
4         *ptr = new;
5     return actual;
6 }
```

**1-core VM:  
Count is correct, no deadlock**

- Spin lock

```

1 int lock(int *ptr) {
2     while (CompareAndSwap(ptr, 1, 1) != 1)
3         ; // spin
4 }
```

- X86 provides "cmpxchg1" compare-and-exchange instruction
  - cmpxchg8b
  - cmpxchg16b

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.89

## TWO MORE "LOCK BUILDING" CPU INSTRUCTIONS

- Cooperative instructions used together to support synchronization on RISC systems
- No support on x86 processors
  - Supported by RISC: Alpha, PowerPC, ARM
- Load-linked (LL)
  - Loads value into register
  - Same as typical load
  - Used as a mechanism to track competition
- Store-conditional (SC)
  - Performs "mutually exclusive" store
  - Allows only one thread to store value

October 15, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L6.90

### LL/SC LOCK

```
1 int LoadLinked(int *ptr) {
2     return *ptr;
3 }
4
5 int StoreConditional(int *ptr, int value) {
6     if (no one has updated *ptr since the LoadLinked to this address) {
7         *ptr = value;
8         return 1; // success!
9     } else {
10        return 0; // failed to update
11    }
12 }
```

- LL instruction loads pointer value (ptr)
- SC only stores if the load link pointer has not changed
- Requires HW support
  - C code is psuedo code

October 15, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.91
------------------	---	-------

### LL/SC LOCK - 2

```
1 void lock(lock_t *lock) {
2     while (1) {
3         while (LoadLinked(&lock->flag) == 1)
4             ; // spin until it's zero
5         if (StoreConditional(&lock->flag, 1) == 1)
6             return; // if set-it-to-1 was a success: all done
7                 otherwise: try it all over again
8     }
9 }
10
11 void unlock(lock_t *lock) {
12     lock->flag = 0;
13 }
```

- Two instruction lock

October 15, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L6.92
------------------	---	-------

# QUESTIONS

